

# **BIOINFORMATICS**

**Mikel Egaña Aranguren**

[www.sindominio.net/~pik](http://www.sindominio.net/~pik)

**Manchester University Bioinformatics MSc**

## **1.- INTRODUCCION:**

- **Historia y bases generales de la Bioinformatica: bases de datos y algoritmos.**

## **2.- PROBLEMAS ACTUALES DE LA BIOINFORMATICA:**

- **Informacion vs. conocimiento.**
- **Crecimiento de bases de datos.**
- **Brecha bioinformaticos/biologos.**

## **3.- ALGUNAS DE LAS ULTIMAS TENDENCIAS EN BIOINFORMATICA :**

- **Hacia una computacion semantica: ontologias y la web semantica.**
- **Hacia una computacion distribuida: myGRID.**
- **Hacia una computacion estandar: los estandares en la investigacion con microarrays.**

# **1.- INTRODUCCION:**

**Attwood, T.K., Miller, C.J. “Bioinformatics goes back to the future”  
(2003). Nature Reviews. 4; 157-161.**

**<http://www.sindominio.net/~pik/bioinformatics.html>**

**Biologia molecular en los 80: secuenciacion, estructuras.**

**Centros principales:**

- **Welcome Trust Genome Campus: EBI (European Bioinformatics Institute), Sanger Institute, Human Genome Mapping Project Resource Center.**
- **National Center for Biological Information (NCBI).**

**Anotacion.**

[ExPASy Home page](#)[Site Map](#)[Search ExPASy](#)[Contact us](#)[Swiss-Prot](#)[Hosted by APAF Australia](#)

Mirror sites:

[Bolivia](#)[Canada](#)[China](#)[Korea](#)[Switzerland](#)[Taiwan](#)[USA](#)

Search

Swiss-Prot/TrEMBL



for

Go

Clear

## NiceProt View of Swiss-Prot: [Q01594](#)

[Printer-friendly view](#)[Submit update](#)[Quick BlastP search](#)
[\[Entry info\]](#)
[\[Name and origin\]](#)
[\[References\]](#)
[\[Comments\]](#)
[\[Cross-references\]](#)
[\[Keywords\]](#)
[\[Features\]](#)
[\[Sequence\]](#)
[\[Tools\]](#)

Note: most headings are clickable, even if they don't appear as links. They link to the [user manual](#) or [other documents](#).

### Entry information

Entry name	ALL1_ALLSA
Primary accession number	Q01594
Secondary accession numbers	None
Entered in Swiss-Prot in	Release 26, July 1993
Sequence was last modified in	Release 26, July 1993
Annotations were last modified in	Release 43, March 2004

### Name and origin of the protein

Protein name	Alliin lyase 1 [Precursor]
Synonyms	EC <a href="#">4.4.1.4</a> Alliinase 1 Cysteine sulphoxide lyase 1
Gene name	None
From	<a href="#">Allium sativum (Garlic)</a> [TaxID: <a href="#">4682</a> ]
Taxonomy	<a href="#">Eukaryota</a> ; <a href="#">Viridiplantae</a> ; <a href="#">Streptophyta</a> ; <a href="#">Embryophyta</a> ; <a href="#">Tracheophyta</a> ; <a href="#">Spermatophyta</a> ; <a href="#">Magnoliophyta</a> ; <a href="#">Liliopsida</a> ; <a href="#">Asparagales</a> ; <a href="#">Alliaceae</a> ; <a href="#">Allium</a>

## References

- [1] SEQUENCE FROM NUCLEIC ACID, AND SEQUENCE OF 39-58.  
TISSUE=[Shoot](#);  
MEDLINE=93049322; PubMed=1385120; [[NCBI](#), [ExpASY](#), [EBI](#), [Israel](#), [Japan](#)]  
[van Damme E., J.M.](#), [Smeets K.](#), [Torrekens S.](#), [van Leuven F.](#), [Peumans W., J.](#);  
"Isolation and characterization of alliinase cDNA clones from garlic (*Allium sativum* L.) and related species.";  
*Eur. J. Biochem.* 209:751-757(1992).
- [2] CHARACTERIZATION, AND CRYSTALLIZATION.  
TISSUE=[Bulb](#);  
MEDLINE=22047089; PubMed=12051663; [[NCBI](#), [ExpASY](#), [EBI](#), [Israel](#), [Japan](#)]  
[Kuettner E.B.](#), [Hilgenfeld R.](#), [Weiss M.S.](#);  
"Purification, characterization, and crystallization of alliinase from garlic.";  
*Arch. Biochem. Biophys.* 402:192-200(2002).
- [3] X-RAY CRYSTALLOGRAPHY (1.53 ANGSTROMS) OF 39-465, AND SUBUNIT.  
TISSUE=[Bulb](#);  
MEDLINE=22336365; PubMed=12235163; [[NCBI](#), [ExpASY](#), [EBI](#), [Israel](#), [Japan](#)]  
[Kuettner E.B.](#), [Hilgenfeld R.](#), [Weiss M.S.](#);  
"The active principle of garlic at atomic resolution.";  
*J. Biol. Chem.* 277:46402-46407(2002).

## Comments

- ◆ **FUNCTION:** Able to cleave the C-S bond of sulfoxide derivatives of Cys to produce allicin, thus giving rise to all sulfur compounds which are responsible for most of the properties of garlic, such as the specific smell and flavor as well as the health benefits like blood lipid or blood pressure lowering.
- ◆ **CATALYTIC ACTIVITY:** An S-alkyl-L-cysteine S-oxide = an alkyl sulfenate + 2 aminocrylate.
- ◆ **COFACTOR:** Pyridoxal phosphate.
- ◆ **SUBUNIT:** Homodimer.
- ◆ **SUBCELLULAR LOCATION:** Vacuolar.
- ◆ **TISSUE SPECIFICITY:** Bulb and shoots.
- ◆ **DOMAIN:** The 6 Cys residues of the EGF-like domain are arranged in a disulfide pattern different from the one found in the canonical EGFs. The

File Edit View Web Go Bookmarks Tabs Help

New Back Forward Stop Refresh Home Fullscreen Zoom 80 http://au.expasy.org/cgi-bin/niceprot.pl?Q01594

Mikel Egaña A x PLoS Biology: x Bioinformatics x EMBER: Chap x NiceProt View x



### Cross-references

EMBL	Z12622; CAA78268.1; -. <a href="#">[EMBL / GenBank / DDB.]</a> <a href="#">[CoDingSequence]</a>
PIR	<a href="#">S29302</a> ; S29302.
PDB	1LK9; 11-DEC-02. <a href="#">[ExPASy / RCSB / EBI]</a>
InterPro	<a href="#">IPR006948</a> ; Alliinase_C. <a href="#">IPR006947</a> ; EGF_alliinase. <a href="#">Graphical view of domain structure.</a>
Pfam	<a href="#">PF04864</a> ; Alliinase_C; 1. <a href="#">PF04863</a> ; EGF_alliinase; 1. <a href="#">Pfam graphical view of domain structure.</a>
PROSITE	<a href="#">PS00022</a> ; EGF_1; 1. <a href="#">PS01186</a> ; EGF_2; FALSE_NEG.
ProDom	<a href="#">[Domain structure / List of seq. sharing at least 1 domain]</a>
BLOCKS	<a href="#">Q01594</a> .
ProtoNet	<a href="#">Q01594</a> .
ProtoMap	<a href="#">Q01594</a> .
PRESAGE	<a href="#">Q01594</a> .
DIP	<a href="#">Q01594</a> .
ModBase	<a href="#">Q01594</a> .
SMR	<a href="#">Q01594</a> ; D7862E867AD74383.
SWISS-2DPAGE	<a href="#">Get region on 2D PAGE.</a>
UniRef	View cluster of proteins with at least <a href="#">50%</a> / <a href="#">90%</a> identity.

### Keywords

[Lyase](#); [Pyridoxal phosphate](#); [Chloride](#); [Signal](#); [EGF-like domain](#); [Glycoprotein](#); [3D-structure](#).

### Features

 [Feature table viewer](#)  [Feature aligner](#)

Done.

**Features**



[Feature table viewer](#)



[Feature aligner](#)

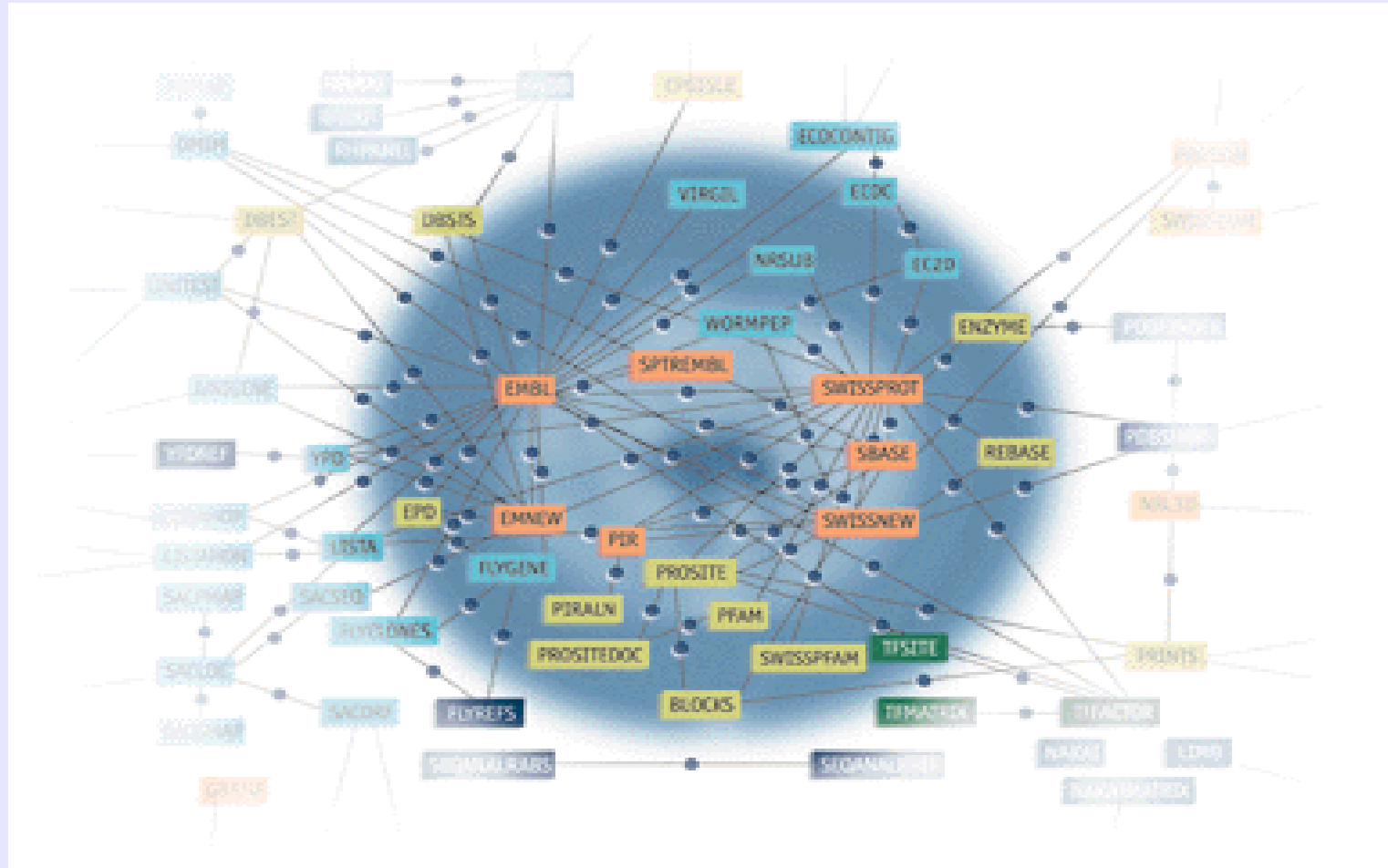
Key	From	To	Length	Description
SIGNAL	<a href="#">1</a>	<a href="#">28</a>	28	Potential.
PROPEP	<a href="#">29</a>	<a href="#">38</a>	10	
CHAIN	<a href="#">39</a>	<a href="#">486</a>	448	Alliin lyase 1.
DOMAIN	<a href="#">51</a>	<a href="#">97</a>	47	EGF-like (atypical).
BINDING	<a href="#">132</a>	<a href="#">132</a>		Chloride.
BINDING	<a href="#">136</a>	<a href="#">136</a>		Chloride.
BINDING	<a href="#">138</a>	<a href="#">138</a>		Chloride.
BINDING	<a href="#">289</a>	<a href="#">289</a>		Pyridoxal phosphate.
DISULFID	<a href="#">58</a>	<a href="#">77</a>		
DISULFID	<a href="#">79</a>	<a href="#">88</a>		
DISULFID	<a href="#">82</a>	<a href="#">95</a>		
DISULFID	<a href="#">406</a>	<a href="#">414</a>		
CARBOHYD	<a href="#">57</a>	<a href="#">57</a>		N-linked (GlcNAc...) (Potential).
CARBOHYD	<a href="#">184</a>	<a href="#">184</a>		N-linked (GlcNAc...).
CARBOHYD	<a href="#">229</a>	<a href="#">229</a>		N-linked (GlcNAc...) (Potential).
CARBOHYD	<a href="#">366</a>	<a href="#">366</a>		N-linked (GlcNAc...).

**Sequence information**

<b>Length: 486 AA [This is the length of the unprocessed precursor]</b>	<b>Molecular weight: 55638 Da [This is the MW of the unprocessed precursor]</b>	<b>CRC64: D7862E867AD74383 [This is a checksum on the sequence]</b>			
10	20	30	40	50	60
MVESYKKIGS	CNKMPCLVIL	TCIIMSNSLV	NNMMVQAKH	TWTMKAAEEA	EAVANINCSE
70	80	90	100	110	120

## Bases de datos:

- Repositorios de secuencias: swiss-prot, trEMBL, GenBank.
- Sistemas que acceden a varias bd-s a la vez: SRS (EBI), Entrez (NCBI).







Search across databases

Kinase

GO CLEAR Help

305015		<b>PubMed:</b> biomedical literature citations and abstracts	?	1573		<b>Books:</b> online books	?
43667		<b>PubMed Central:</b> free, full text journal articles	?	1804		<b>OMIM:</b> Online Mendelian Inheritance in Man	?
				227		<b>Site Search:</b> NCBI web and FTP sites	?
828565		<b>Nucleotide:</b> sequence database (GenBank)	?	12539		<b>UniGene:</b> gene-oriented clusters of transcript sequences	?
81288		<b>Protein:</b> sequence database	?	312		<b>CDD:</b> conserved protein domain database	?
383		<b>Genome:</b> whole genome sequences	?	4122		<b>3D Domains:</b> domains from Entrez Structure	?
1003		<b>Structure:</b> three-dimensional macromolecular structures	?	4748		<b>UniSTS:</b> markers and mapping data	?
none		<b>Taxonomy:</b> organisms in GenBank	?	82		<b>PopSet:</b> population study data sets	?
234422		<b>SNP:</b> single nucleotide polymorphism	?	214530		<b>GEO:</b> expression and molecular abundance profiles	?
18309		<b>Gene:</b> gene-centered information	?	13		<b>GEO DataSets:</b> experimental sets of GEO data	?
12442		<b>HomoloGene:</b> Eukaryotic homology groups	?	40		<b>Cancer Chromosomes:</b> cytogenetic databases	?
none		<b>Journals:</b> detailed information about the journals indexed in PubMed and other Entrez databases	?	84		<b>MeSH:</b> detailed information about NLM's controlled vocabulary	?

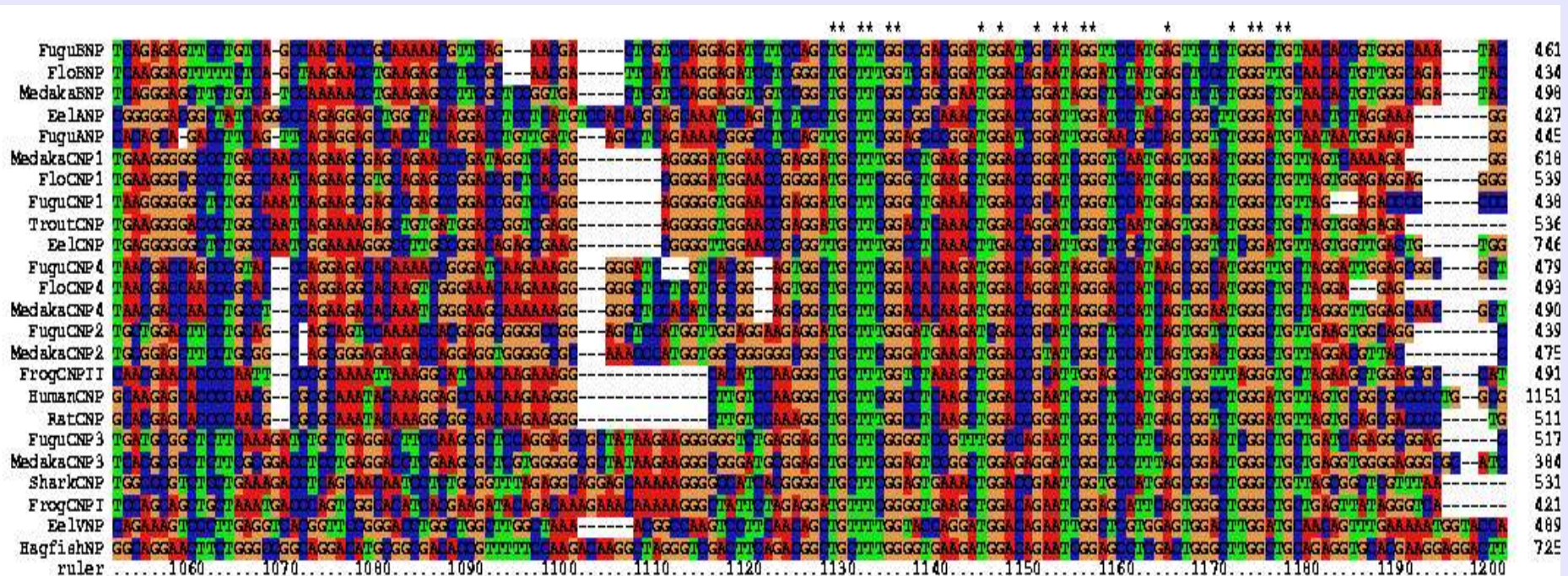
## **Bases de datos de familias, patrones o motifs: abstraccion.**

- **Mejor rendimiento en bajas similitudes.**
- **Busquedas funcionales.**
- **Unicos motifs: prosite, emotif. Expresiones regulares:  
D-M-x-[ILV]-x{2}-G**
- **Dominio: profiles (profiles), pfam (hidden markov models).**
- **Multiples motifs: fingerprints, blocks.**
- **Interpro.**

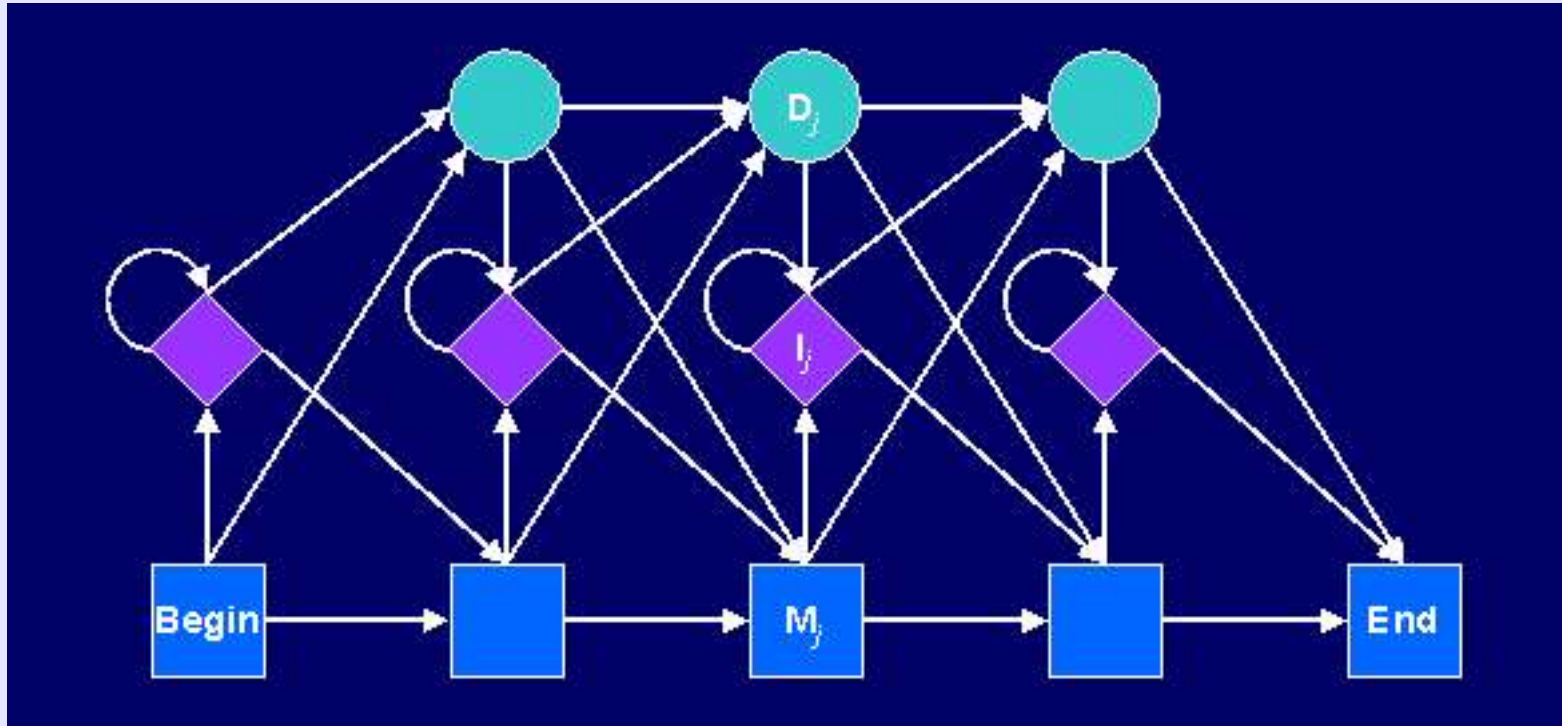
**Bases de datos de estructuras, etc.**

# Algoritmos:

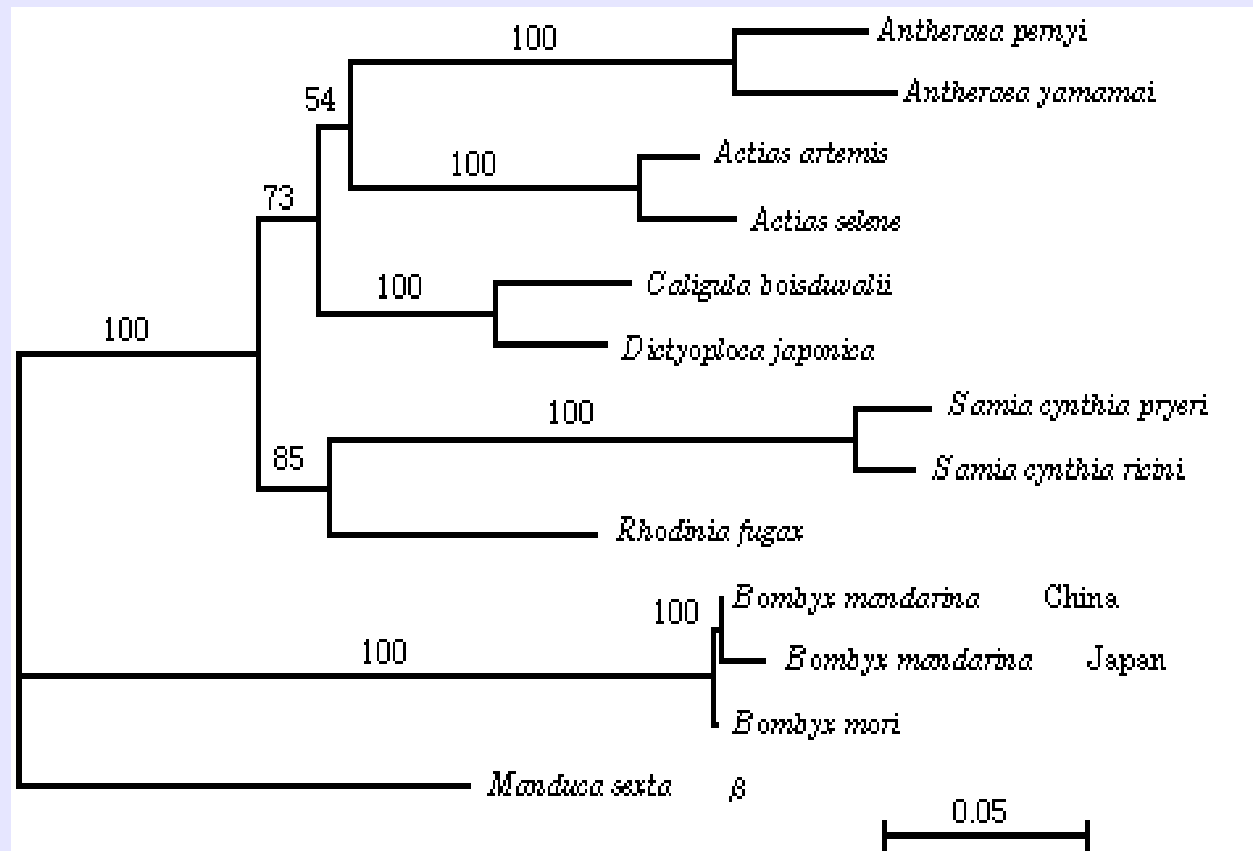
- Prediccion de estructuras: CASP5 (Critical Assesment of technics for protein Structure Prediction): <http://www.predictioncenter.llnl.gov/>
- Alineamientos: algoritmos recursivos. Alineamiento local (BLAST: Basic Local Alignment Search Tool) o global. Matrices de similitud.



- **Hidden Markov Models:**



- **Algoritmos filogeneticos:**



- **ETC ...**

- **BioPerl:** <http://www.bioperl.org>.

“How Perl saved the human genome project”:

[http://www.stanford.edu/class/gene211/handouts/How\\_Perl\\_HGP.html](http://www.stanford.edu/class/gene211/handouts/How_Perl_HGP.html)

## **2.- PROBLEMAS:**

### **Informacion vs. conocimiento:**

- **Infoberg: las bd-s duplican su contenido cada 9 meses.**
- **Anotacion: errores multiplicados, revisiones imposibles. Anotacion automatica. Similitud - funcion ???**
- **Homologos: ortologos - paralogos.**
- **Definicion de gen.**
- **Definicion de funcion.**
- **Naturaleza modular de las proteinas.**
- **Redundancia: NRDB: non-identical but redundant.**

**Integracion de recursos: contingencia historica. N scripts son demasiados scripts para N herramientas.**

**Brecha bioinformaticos – biologos:**

- **Falta de cultura computacional:**

[http://sindominio.net/biblioweb/telematica/command\\_es](http://sindominio.net/biblioweb/telematica/command_es)

### 3.-LINEAS DE FUTURO:

Hacia una computacion semantica:

Ontologias: vocabulario + relaciones estructuradas:

Sistematizacion del conocimiento.

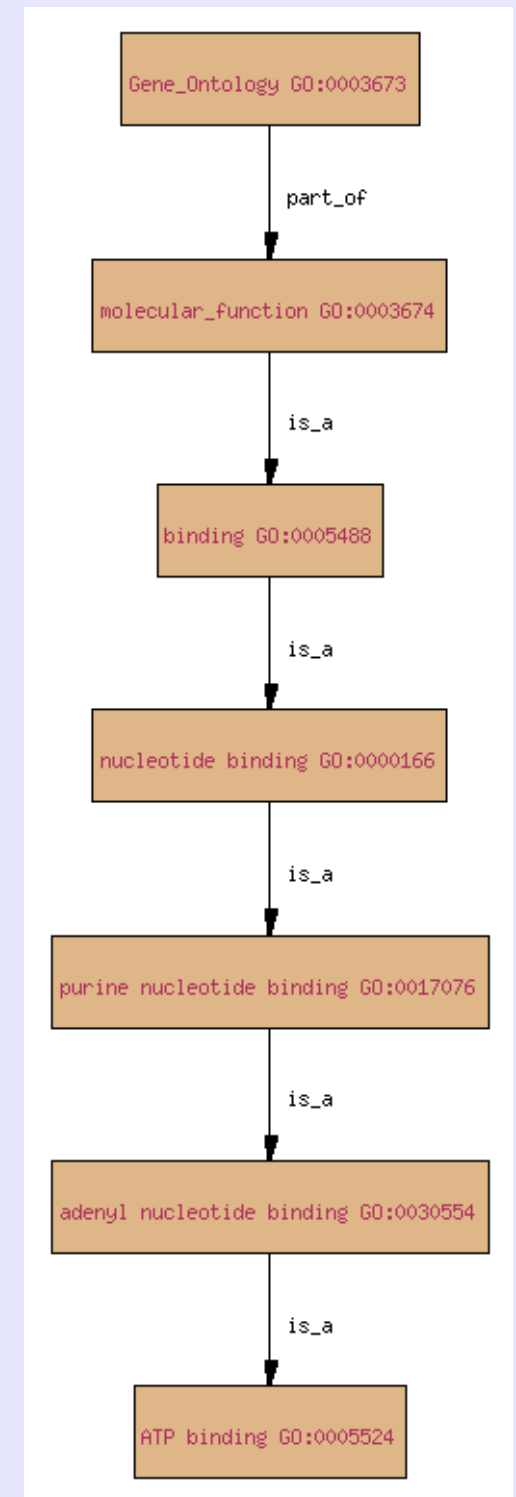
Semantica computacionalmente accesible.

Gene Ontology: proceso biologico, componente celular, funcion molecular.

<http://www.geneontology.org>

Open Biological Ontologies:

<http://obo.sourceforge.net/>





## La web semantica:

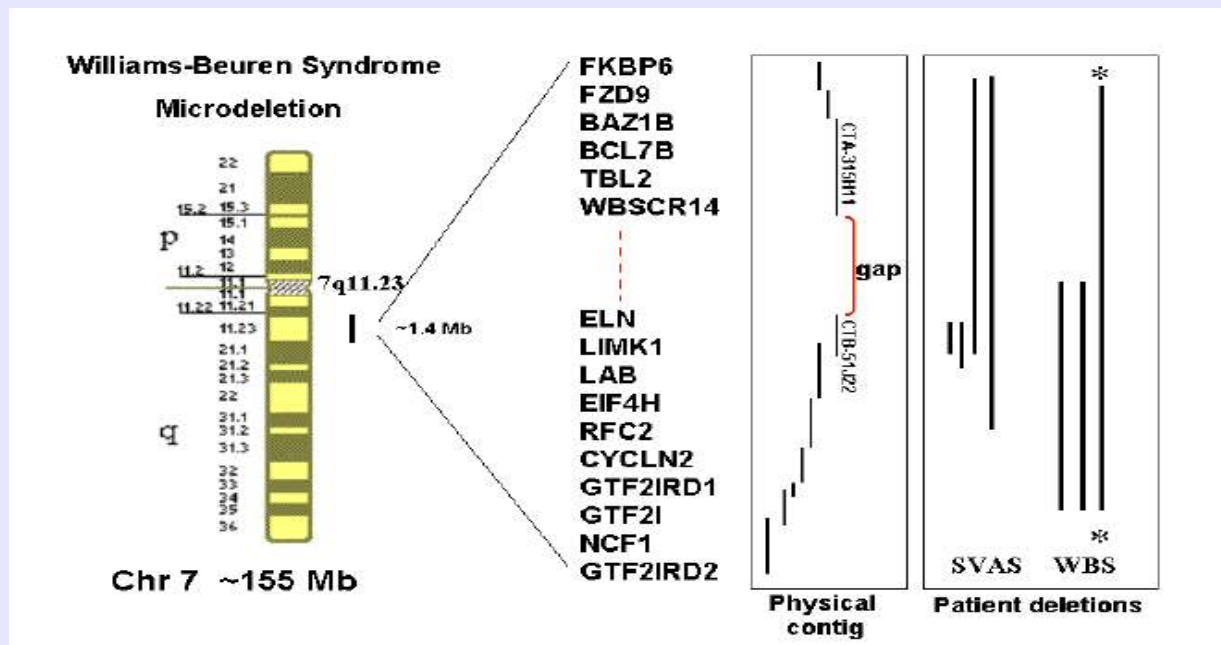
- **W3C: <http://www.w3.org/>. World Wide Web Consortium.**
- **"The Semantic Web is an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation." -- Tim Berners-Lee, James Hendler, Ora Lassila, The Semantic Web, Scientific American, May 2001.**
- **RDF (Resource Description Framework) + XML (eXtensible Markup Language).**
- **OWL (Web Ontology Language).**

## **Hacia una computacion distribuida:**

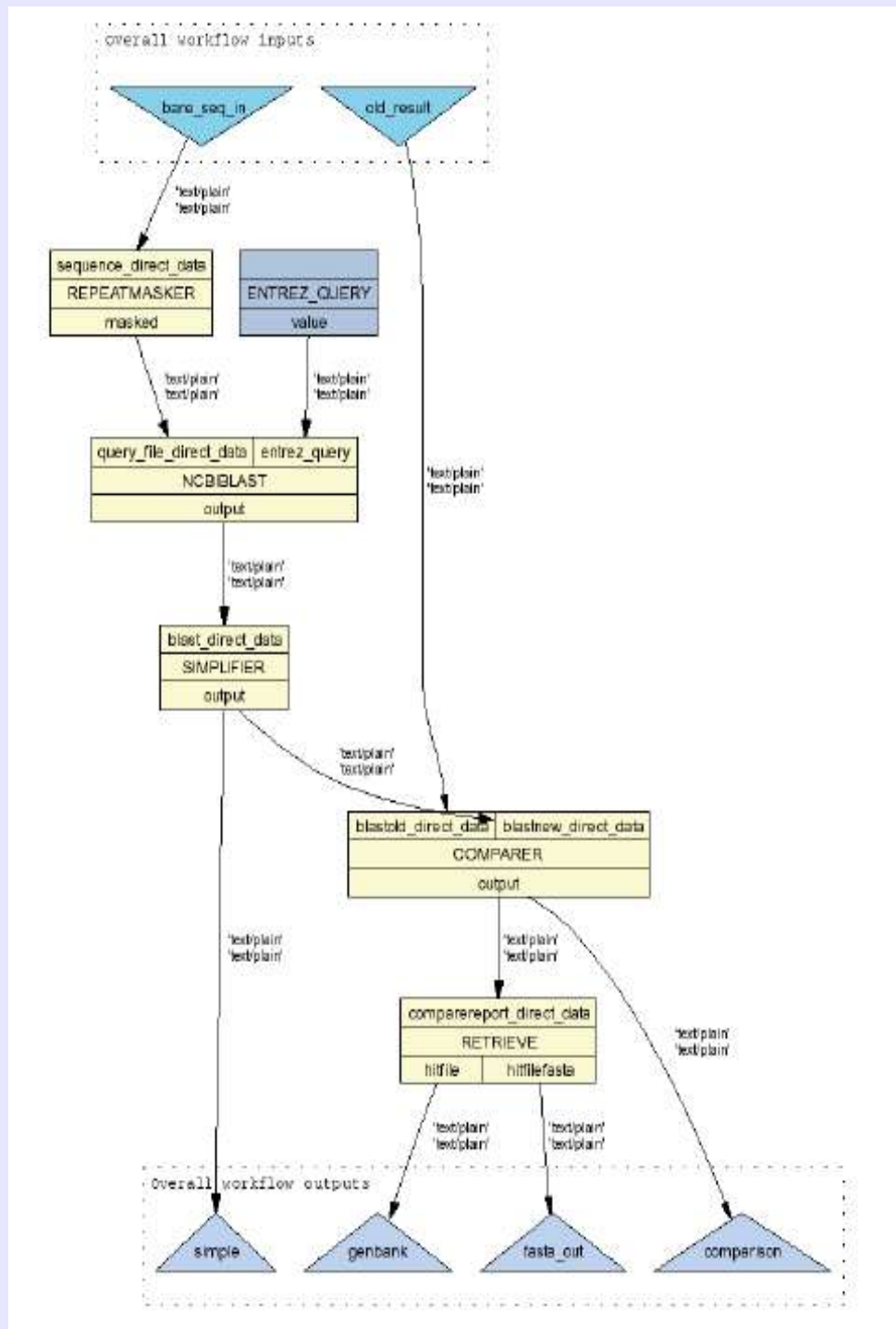
- **Grid: computacion distribuida entre diferentes centros (con diferentes lenguajes de programacion, plataformas, ...) con una serie de protocolos estandar.**
- **MyGrid: entorno virtual, automatizado y distribuido de experimentos in silico. <http://www.mygrid.org.uk/>**
- **Entorno comun de investigacion, automatico (rapidez), almacenaje de datos, distribuido y semantico, ...**
- **R. Stevens, H.J. Tipney, C. Wroe, T. Oinn, M. Senger, P. Lord, C.A. Goble, A. Brass and M. Tassabehji “Exploring Williams-Beuren Syndrome Using myGrid” to appear in Proceedings of 12th International Conference on Intelligent Systems in Molecular Biology, 31st Jul-4th Aug 2004, Glasgow, UK.**

## Sindrome Williams – Beuren:

- Deleciones cromosoma 7.
- Zona altamente repetitiva y compleja: WBS critical region, todavia sin caracterizar.
- Analisis periodicos >>> secuencias nuevas de esa zona >>> bateria de herramientas bioinformaticas >>> almacenaje de datos y procedimientos en el “LabBook”.

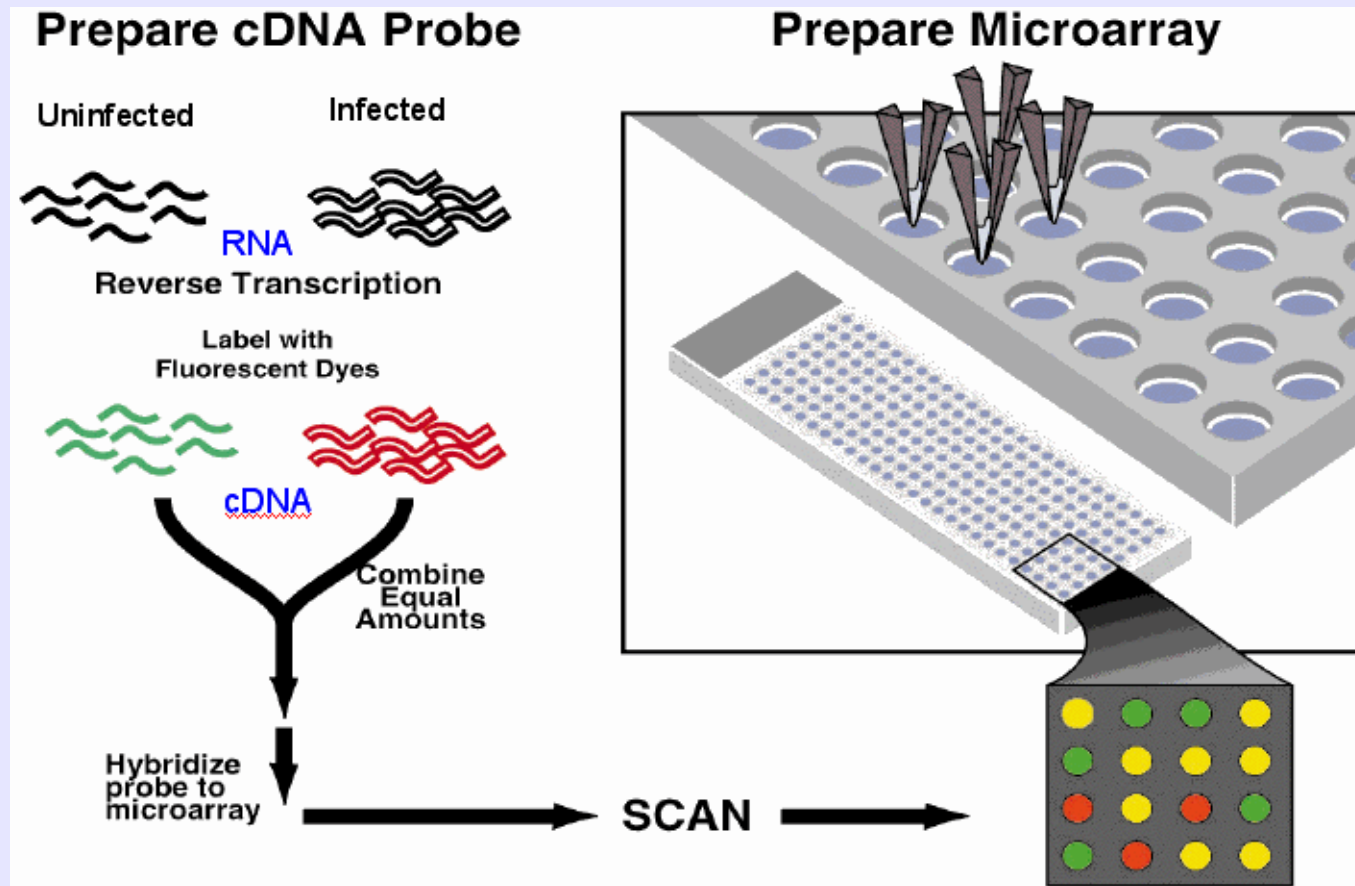


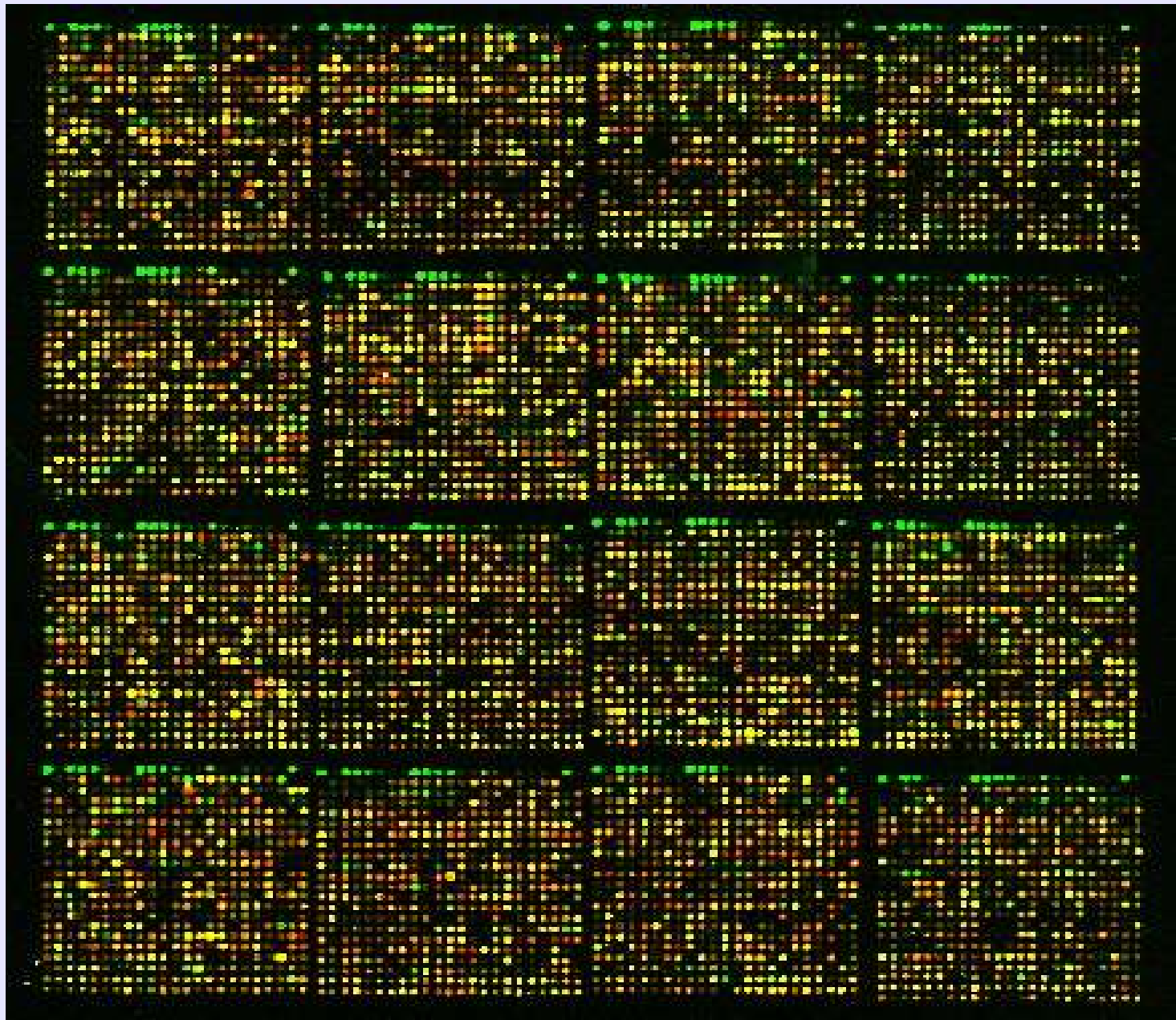
# Workflows



# Hacia una computacion estandar:

- **Microarrays:**





## **Babelismo en bases de datos:**

- **MGED: Microarray Gene Expression Data:** <http://www.mged.org/>.  
Estandares en el analisis, anotacion, almacenamiento y dispersion de datos sobre expresion de genes.
- **MIAME: Minimal Information About a Microarray Experiment.**  
Directivas para la anotacion de los resultados con microarrays.  
Adoptado por Nature, Cell. Doloroso para los biologos, bueno para la comunidad.
- **MAGE-ML: Microarray Gene Expression Markup Language.**
- **MAGE-OM: Microarray Gene Expression Object Model.**